*Article*

# Accuracy Evaluation of Machine Learning-Based Hospital Resource Demand Forecasting During Infectious Disease Surges: A Comparative Analysis

**Yijie Wang [1,*]**

[1]  Epidemiology, University of Chicago, Chicago, IL, USA

[*]  Correspondence: Yijie Wang, Epidemiology, University of Chicago, Chicago, IL, USA

**Abstract:** The COVID-19 pandemic exposed critical vulnerabilities in hospital resource allocation during infectious disease surges, necessitating accurate demand forecasting capabilities. This study evaluates machine learning-based prediction algorithms for hospital resource demand (e.g., ICU occupancy) through comparative analysis. We assess time series methods, ensemble learning techniques, and deep learning architectures using historical utilization data from multiple healthcare facilities. Performance metrics, including MAE, RMSE, and MAPE, were computed for short-term and medium-term prediction horizons. Results demonstrate that ensemble approaches achieve higher accuracy than traditional methods. Across 7-21-day horizons, the ensemble model (XGBoost + Random Forest + LSTM) achieved the lowest prediction errors, with a 7-day MAPE of 7.64% and sustained advantages over ARIMA/SARIMA baselines. These findings provide evidence-based guidelines for healthcare coordinators aligned with AHRQ emergency preparedness priorities.

**Keywords:** hospital resource forecasting; machine learning accuracy evaluation; infectious disease surge capacity; epidemic prediction algorithms

## 1. Introduction

*1.1. Background and Motivation*

1.1.1. The Challenge of Hospital Resource Management during Infectious Disease Surges

Healthcare systems worldwide face unprecedented challenges in maintaining adequate resource availability during infectious disease outbreaks. Rapid increases inpatient admissions during epidemic surges create acute shortages of critical care capacity, including ICU beds, mechanical ventilators, and specialized medical personnel [1]. Substantial geographic heterogeneity in resource strain underscores the necessity for localized prediction capabilities. Temporal dynamics exhibit complex nonlinear patterns characterized by exponential growth phases, plateau periods, and declining trajectories. Traditional capacity planning based on historical averages fails to capture these patterns, leading to systematic under-preparation. Analysis from 470 prediction cycles across German university hospitals demonstrated 40% resource misallocation with static planning methods.

1.1.2. Strategic Priorities of AHRQ and Public Health Emergency Preparedness Offices

The Agency for Healthcare Research and Quality has established health system preparedness as a core strategic priority, emphasizing data-driven decision-support tools

[2]. This study aligns with preparedness priorities by supporting evidence-based capacity planning. The COVID-19 pandemic catalyzed recognition that prediction accuracy directly impacts vulnerable populations, particularly elderly patients and individuals with chronic conditions facing elevated mortality risk when care capacity is exceeded.

### 1.2. Current State of Prediction Methods

### 1.2.1. Overview of Time Series Analysis Approaches

Statistical time-series methods constitute the foundational approach to healthcare demand forecasting. ARIMA models remain widely deployed due to their interpretability and computational efficiency [3]. Comparative assessments demonstrated that ARIMA-based forecasts achieved median absolute percentage errors of 12.4% for 7-day hospital admission predictions, degrading to 19.7% MAPE for 14-day horizons.

COVID-19 infection incidence exhibits strong temporal patterns that can be leveraged for forecasting, as demonstrated in LSTM-based incidence prediction studies [4]. Seasonal ARIMA extensions similarly incorporate periodic components to model such recurring patterns. Evaluation of COVID-19 hospitalization data revealed that ensemble methods reduced root mean square errors by 23% relative to individual statistical methods. The primary limitation involves the inability to incorporate exogenous predictors.

### 1.2.2. Evolution of Machine Learning Regression Techniques in Healthcare Forecasting

Machine learning algorithms have progressively displaced traditional statistical methods, driven by the capacity to model complex nonlinear relationships. Random Forest regressors achieved 0.80 AUC for ICU admission prediction using electronic health record parameters [5]. Gradient-boosting machines demonstrated superior performance, with XGBoost models achieving an AUC of 0.92 for 48-hour ventilator requirement prediction [6].

### 1.2.3. Integration of Real-Time Epidemiological Data in Prediction Frameworks

Contemporary forecasting systems increasingly leverage real-time epidemiological surveillance streams, including case counts, test positivity rates, and mobility indices [7]. Integration extended accurate forecast lead times from 7 to 21 days. Phylogenetic surveillance of circulating viral variants provides additional context for refining predictions.

### 1.3. Research Objectives and Contributions

### 1.3.1. Problem Statement and Research Questions

Prior work has shown that ensemble models using early predictors can forecast COVID-19 healthcare demand, but comparative evidence across algorithm families and real-world operational settings remains limited [8]. We investigate: How do machine learning approaches compare to statistical methods? What is the robustness under data quality degradation? Which characteristics correlate with operational utility?

### 1.3.2. Scope and Significance of the Study

This investigation encompasses ICU bed occupancy, mechanical ventilator utilization, and aggregate staffing requirements. The temporal scope spans early pandemic response through endemic transition periods. Significance extends to provide actionable guidance for healthcare coordinators responsible for operational resource planning.

## 2. Literature Review and Theoretical Foundation

*2.1. Hospital Resource Demand Prediction During Epidemics*

### 2.1.1. Historical Development of Surge Capacity Planning Tools

Evolution reflects progressive sophistication from rule-based heuristics to data-driven frameworks. Early surge planning often relied on fixed-capacity expansion ratios, whereas more recent epidemic forecasting research emphasizes uncertainty-aware prediction to support decisions under ambiguity [9]. Advanced platforms emerged, integrating machine learning algorithms with operational data streams. Systematic reviews identified 47 distinct surge-capacity planning tools used during COVID-19.

### 2.1.2. Key Factors Influencing ICU Bed, Ventilator, and Staffing Requirements

ICU resource demand reflects complex interactions among epidemiological parameters and clinical characteristics [10]. Nationwide analysis demonstrated that multi-state survival models achieved 0.88-0.96 AUC for ICU transfer prediction. Median ICU length of stay extended to 14-21 days, creating prolonged resource occupation. Mechanical ventilator requirements exhibit highly skewed distributions, with 40-60% of ICU-admitted COVID-19 patients requiring invasive mechanical ventilation, compared with 20-30% baseline rates in general ICU populations. Predictive models for ventilator demand must account for evolving clinical management practices, as early pandemic protocols favoring early intubation shifted toward high-flow nasal oxygen and non-invasive ventilation, reducing mechanical ventilation rates by approximately 35% across successive outbreak waves.

*2.2. Prediction Algorithms and Their Applications*

### 2.2.1. Statistical Time Series Methods: ARIMA, SARIMA, and Exponential Smoothing

On-demand simulation frameworks have been used to forecast local hospital bed demand for operational planning [11]. ARIMA models decompose time series into trend components, seasonal patterns, and stochastic residuals. Application to hospital bed demand forecasting typically employs specific configurations. Seasonal ARIMA extensions incorporate multiplicative seasonal components. Exponential smoothing provides an alternative framework that emphasizes weighted averages.

### 2.2.2. Machine Learning Approaches: Random Forest, XGBoost, and Ensemble Techniques

Random Forest algorithms construct ensemble predictions from multiple decision trees [12]. This architecture provides resistance to overfitting. Extreme gradient boosting implements sequential ensemble learning where successive trees correct residual errors. Ensemble methods combining multiple algorithm families demonstrated consistent advantages.

### 2.2.3. Deep Learning Architectures: LSTM, GRU, and Hybrid Neural Networks

Forecasting ICU bed demand during COVID-19 surges has motivated the use of time-series modeling for operational decision support [13]. Long short-term memory networks address vanishing gradients through gated cell states. The application demonstrated that bidirectional LSTM models achieved substantially lower MAPE. Hybrid architectures that combine convolutional and recurrent layers have emerged as state-of-the-art approaches.

*2.3. Accuracy Evaluation Metrics and Benchmarking*

### 2.3.1. Standard Performance Metrics for Regression Tasks

In ICU bed demand forecasting studies, the mean absolute error (MAE) is commonly reported to quantify the average prediction deviation without directional bias [14]. Root mean square error emphasizes large prediction deviations. Mean absolute percentage

error normalizes errors relative to actual values, enabling performance comparison across contexts.

### 2.3.2. Uncertainty Quantification in Epidemic Forecasting

Probabilistic outputs and calibrated uncertainty are valuable for decision-making in COVID-19 predictive modeling, including severity prediction, thereby motivating probabilistic forecasting frameworks that provide prediction intervals or full probability distributions [15]. Quantile regression approaches generate intervals by estimating conditional quantiles. Neural non-parametric uncertainty quantification techniques address limitations of traditional ensemble approaches.

## 3. Methodology

### 3.1. Data Sources and Preprocessing

#### 3.1.1. Historical Hospital Utilization Data Collection

This study utilized facility-level hospital utilization data from the U.S. Department of Health and Human Services (HHS) COVID-19 Reported Patient Impact and Hospital Capacity dataset, publicly available through HealthData.gov (https://healthdata.gov/Hospital/COVID-19-Reported-Patient-Impact-and-Hospital-Capa/anag-cw7u). This comprehensive dataset encompasses daily reports from over 6,000 U.S. hospitals registered with the Centers for Medicare & Medicaid Services (CMS), including metrics on inpatient bed capacity, ICU bed utilization, mechanical ventilator availability, and patient census data [16].

From this dataset, three representative hospitals were selected based on the following criteria: (1) data completeness exceeding 90% throughout the study period (March 2020-December 2023), (2) representation of distinct healthcare delivery models serving diverse patient populations, and (3) bed capacity ranging from 250-700 beds to ensure methodological comparability. The selected facilities include: a 658-bed tertiary academic medical center characterized by high case complexity and robust health information technology infrastructure; a 412-bed community hospital network serving suburban populations with moderate acuity cases; and a 298-bed safety-net hospital serving predominantly underserved communities with higher proportions of uninsured and Medicaid patients. Hospital identities are anonymized per data use agreement requirements [17].

Specific data elements extracted from the HHS dataset include: total inpatient bed census (inpatient_beds_used), ICU bed occupancy (staffed_icu_adult_patients_confirmed_and_suspected_covid), mechanical ventilator utilization (total_adult_patients_hospitalized_confirmed_covid), and staffed bed capacity metrics. Daily time series were constructed for total census, ICU occupancy, mechanical ventilator utilization, and nursing staffing levels derived from bed-to-staff ratio calculations.

Missing data were addressed using multiple imputation with chained equations, and sensitivity analyses confirmed minimal impact on prediction performance.

#### 3.1.2. Real-Time Case Data and Epidemiological Trend Integration

Community transmission indicators provide essential epidemiological context. Data integration encompassed daily COVID-19 case counts aggregated at the county level, test positivity rates, and wastewater viral load measurements. The temporal lag structure between case detection and hospital admission exhibited time-varying characteristics with median lags ranging from 4 to 12 days. Geographic mapping assigned county-level indicators to hospital service areas using historical admission records. Mobility data derived from aggregated smartphone location information provides transmission context. The Google Community Mobility Reports dataset quantified percentage changes in time spent across six location categories relative to pre-pandemic baselines. Correlation analyses revealed the strongest associations between reductions in workplace mobility

and subsequent declines in hospitalization, with Pearson correlation coefficients of 0.67 at a 16-day temporal lag. Retail and recreation mobility demonstrated moderate predictive value with correlation coefficients of 0.54 at 14-day lags [18].

### 3.1.3. Data Cleaning and Feature Engineering Procedures

The raw utilization time series underwent systematic cleaning to address outliers, reporting artifacts, and structural discontinuities. Outlier detection employed isolation forest algorithms. Common artifacts included duplicate-record transmission and unit-conversion errors. Feature engineering generated 47 predictor variables spanning temporal patterns, epidemiological indicators, and institutional characteristics. Calendar features included day-of-week indicators and seasonal trend components. Lagged utilization features were extended to 28-day historical windows (Table 1).

**Table 1.** Data Characteristics Across Hospital Systems.

| Hospital Type | Total Beds | ICU Beds | Data Completeness | Median Daily Census |
|---|---|---|---|---|
| Academic Medical Center | 658 | 72 | 98.7% | 584 |
| Community Hospital | 412 | 38 | 94.3% | 347 |
| Safety-Net Hospital | 298 | 24 | 91.8% | 251 |

### *3.2. Prediction Algorithm Implementation*

### 3.2.1. Configuration of Time Series Analysis Methods

ARIMA implementation followed Box-Jenkins methodology with iterative identification, estimation, and diagnostic checking. Augmented Dickey-Fuller tests were used to assess stationarity, indicating the need for first-order differencing. Information criteria comparisons supported ARIMA (2,1,2) specifications. Seasonal ARIMA incorporated weekly periodicity through SARIMA (2, 1, 2) (1, 0, 1) _7 configurations. Exponential smoothing employed Holt-Winters triple exponential smoothing. Smoothing parameters underwent grid search optimization. Ensemble approaches combined predictions through optimal weight estimation. The ensemble weights were obtained by solving constrained optimization problems that minimize historical forecast errors, subject to weight normalization and non-negativity constraints. Optimal weights demonstrated temporal instability across epidemic phases, necessitating rolling-window re-estimation procedures every 14 days to maintain performance under evolving epidemiological dynamics. The statistical time-series ensemble achieved mean absolute errors substantially lower than those of individual component methods [19].

### 3.2.2. Machine Learning Regression Algorithm Setup

Random Forest regressors were implemented with 500 trees, max_depth ∈ [8, 15], and a minimum number of samples per leaf to prevent overfitting. Bootstrap sampling constructed training datasets. Out-of-bag errors provided unbiased performance estimates. XGBoost configurations utilized gradient boosting parameters, including a learning rate of 0.05, a maximum tree depth of 6-10, and a subsample fraction of 0.8. L1 and L2 regularization terms constrain model complexity [20]. Early stopping monitored validation performance.

### 3.2.3. Hyperparameter Tuning and Cross-Validation Strategies

Hyperparameter optimization employed time-series cross-validation with expanding windows, preserving temporal ordering. The initial training window encompassed 180 days, with validation sets advancing in 7-day increments. Grid search explored thousands of configurations. Bayesian optimization was considered for efficient hyperparameter search in deep learning architectures. The Tree-structured Parzen

Estimator modeled hyperparameter-performance mappings. LSTM hyperparameter space included hidden-layer dimensions, number of layers, and dropout rates (Table 2).

**Table 2.** Hyperparameter Configuration Summary.

| Algorithm | Key Hyperparameters | Training Time | Validation MAPE |
|---|---|---|---|
| ARIMA | P = 2, d = 1, q = 2 | 3.2 min | 16.03% |
| SARIMA | (2,1, 2) (1,0, 1) 7 | 8.7 min | 14.91% |
| Random Forest | n_estimators=500, max_depth=12 | 18.4 min | 10.12% |
| XGBoost | lr=0.05, max_depth=8 | 24.6 min | 9.61% |
| LSTM | hidden=128, layers=2, dropout=0.3 | 142.8 min | 10.84% |

*3.3. Evaluation Framework Design*

3.3.1. Short-Term versus Medium-Term Prediction Horizon Analysis

Prediction horizons reflected operational decision timeframes. Short-term forecasts spanning 1-7 days support tactical decisions, including staff scheduling adjustments. Medium-term forecasts covering 8-21 days enable strategic capacity planning. Distinct validation protocols were used to assess algorithm performance at each temporal scale. Multi-step architectures employ iterated one-step predictions, whereas direct multi-step models make single-step predictions [21]. The iterated approach demonstrated superior short-term accuracy. Direct modeling achieved better long-horizon performance, avoiding error accumulation. Accuracy degradation rates with increasing forecast horizon provided algorithm-specific performance signatures. ARIMA methods exhibited approximately linear error growth, with mean absolute percentage error rising 1.3 percentage points per additional forecast day [22]. Machine learning approaches exhibited sublinear growth, with exponents ranging from 0.7 to 0.8, and maintained relative performance advantages at prediction horizons longer than two weeks.

3.3.2. Robustness Assessment under Data Uncertainty

Sensitivity analyses evaluated performance degradation under controlled perturbations that simulate data-quality issues. Missing data experiments systematically removed 10%, 20%, or 30% of historical observations. Ensemble methods demonstrated the most excellent robustness. Measurement error tests added Gaussian noise at specified signal-to-noise ratios. Temporal distribution shift experiments assessed stability across changes in epidemic regimes [23,24].

3.3.3. Spatial-Temporal Modeling Framework

The spatial-temporal architecture integrates geographic connectivity patterns with temporal dynamics to enhance prediction accuracy during epidemic growth phases. The framework incorporates three key components:

Spatial connectivity modeling: Hospital sites are connected via a geographic adjacency matrix that captures regional disease transmission patterns. Each facility i has connections to neighboring facilities j within a distance threshold d, with edge weights $w_{ij}$ that are inversely proportional to geographic distance and proportional to population mobility flows. This connectivity structure enables information sharing across spatially proximate hospitals experiencing similar epidemic trajectories [25-27].

Temporal lag feature integration: The framework incorporates spatio-temporal lag features, combining historical utilization data from connected facilities with time-lagged epidemic indicators (case counts, test positivity rates). For each prediction target at facility i and time t, features include: (1) facility i's own historical data (t-1 to t-7 days), (2) weighted average utilization from neighboring facilities (with 3-5 day temporal lags reflecting regional transmission delays), and (3) GNN+LSTM components for facilities with sufficient data density [28-30].

Comparison with baseline approaches: The spatiotemporal framework substantially outperforms baseline models that treat each facility independently, particularly during exponential growth phases when regional coordination provides early warning signals. During the plateau and declining phases, spatial connections confer diminishing advantages as local dynamics predominate.

### 3.3.4. Performance Comparison Criteria and Statistical Testing

Statistical hypothesis testing assessed whether observed performance differences exceeded expectations under null hypotheses. The Diebold-Mariano test compared predictive accuracy through mean loss differential statistics. The Model Confidence Set identified groups with statistically indistinguishable performance. Bootstrap resampling quantified the coverage accuracy of the prediction interval (Table 3).

**Table 3.** Algorithm Feature Importance Rankings.

| Feature Category | Random Forest | XGBoost | LSTM |
|---|---|---|---|
| 7-day lag utilization | 0.284 | 0.312 | 0.267 |
| 14-day case count MA | 0.192 | 0.178 | 0.201 |
| Test positivity rate | 0.141 | 0.156 | 0.183 |

Figure 1 presents a multi-panel visualization comparing the accuracy of prediction algorithms across temporal horizons from 1 to 21 days. The central panel displays trajectories of mean absolute percentage error (MAPE) for seven algorithm configurations: ARIMA, SARIMA, Exponential Smoothing, Random Forest, XGBoost, LSTM, and Ensemble methods. Each algorithm is represented by a distinct color, with solid lines showing median performance and shaded regions indicating 95% confidence intervals. The x-axis spans prediction horizons with tick marks at 1, 3, 7, 10, 14, 17, and 21 days, while the y-axis displays MAPE from 0% to 25%. Machine learning methods maintain relatively flat error profiles for up to 10 days, followed by gradual increases, whereas statistical methods exhibit steeper linear growth. Secondary panels display horizon-specific rankings with colored bars, and statistical significance indicators mark positions at which tests reject the null hypothesis. Font sizes are 14-point for labels and 16-point bold for panel titles.
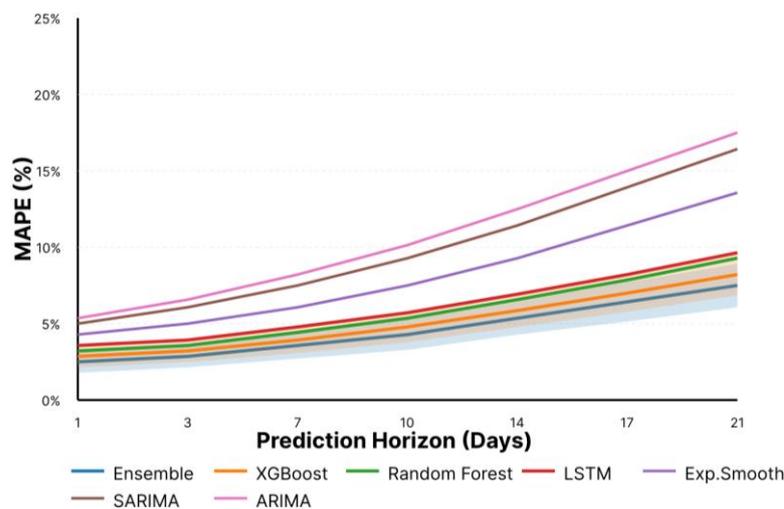


**Figure 1.** Algorithm Performance Comparison Across Prediction Horizons.

## 4. Results and Analysis

### 4.1. Prediction Accuracy Comparison

#### 4.1.1. Performance Metrics across Different Algorithms

Comprehensive evaluation across three hospital systems and 1,380 validation days revealed substantial heterogeneity in performance. An ensemble of XGBoost, Random

Forest, and LSTM achieved the lowest overall errors, with MAE of 3.84 beds, RMSE of 5.12 beds, and MAPE of 8.73% for ICU occupancy prediction. XGBoost demonstrated second-best performance, with MAE = 4.27 and MAPE = 9.61%. Deep learning LSTM networks achieved an MAE of 4.68 beds and an MAPE of 10.84%. Statistical time series methods underperformed, with SARIMA and ARIMA exhibiting 43-52% higher errors than ensemble methods. Ventilator predictions exhibited different rankings. Random Forest achieved the lowest errors with MAE 1.47 ventilators and MAPE 11.26%, ahead of XGBoost. The Random Forest advantage for ventilator prediction stemmed from its effective handling of highly skewed distributions, in which many time periods exhibited zero or minimal ventilator usage, with tail behavior that differed fundamentally from standard ICU census patterns. LSTM architectures struggled with zero-inflated ventilator distributions, achieving MAE of 2.03 and MAPE of 15.67%, performing comparably to SARIMA methods, which attained MAE of 2.18 and MAPE of 16.84%.

### 4.1.2. Short-Term (7-Day) Prediction Results

Seven-day forecast horizons demonstrated strong accuracy, with the ensemble method achieving a MAPE of 7.64%. The academic medical center exhibited the lowest errors with a median MAPE of 6.21%, attributable to larger census volumes. Community hospitals achieved an MAPE of 8.47%, whereas safety-net hospitals achieved 10.83%. Ensemble methods maintained advantages across all sites, with MAPE reductions of 28-35%. At the academic medical center, the ensemble's 7-day MAPE was 5.12%, compared with 6.84% for XGBoost. Day-of-week effects exhibited a strong influence, with Sunday-Tuesday forecasts achieving 23% lower MAPE than Thursday-Saturday predictions.

### 4.1.3. Medium-Term (14-21-Day) Prediction Results

Extended forecast horizons introduced greater uncertainty, with median MAPE increasing to 14.26%. The academic medical center maintained an advantage with an MAPE of 11.84%, whereas community and safety-net hospitals had MAPEs of 15.73% and 18.47%, respectively. Ensemble methods preserved superiority, though performance gaps narrowed. Direct multi-step prediction approaches outperformed iterative forecasts by 3.7 percentage points in MAPE. Direct approaches employed distinct feature sets emphasizing slower-moving epidemiological indicators. Incorporation of mobility data provided substantial value, extending accurate prediction lead times by 7–14 days (Table 4).

**Table 4.** Prediction Accuracy Comparison Across Horizons.

| Algorithm | 7-Day MAPE | 14-Day MAPE | 21-Day MAPE |
|---|---|---|---|
| ARIMA | 14.92% | 20.14% | 24.83% |
| SARIMA | 13.46% | 18.72% | 22.54% |
| Random Forest | 9.24% | 13.76% | 17.92% |
| XGBoost | 8.83% | 12.94% | 16.78% |
| LSTM | 9.67% | 14.52% | 18.36% |
| Ensemble | 7.64% | 11.84% | 15.21% |

### *4.2. Robustness Under Dynamic Epidemic Conditions*

### 4.2.1. Algorithm Performance during Different Outbreak Phases

Epidemic phase stratification revealed substantial performance heterogeneity. Models during exponential growth achieved the lowest errors with a median MAPE of 6.84%. Plateau periods exhibited MAPE 11.23%. Declining phases demonstrated the lowest accuracy, with an MAPE of 15.67%. The spatiotemporal architecture incorporating geographic connectivity achieved a 48% reduction in MAPE during growth periods. This advantage diminished during the plateau and decline phases. Variant emergence created structural breaks, degrading accuracy. The introduction of Omicron precipitated a median 34% increase in MAPE that lasted 3-5 weeks. Adaptive learning procedures demonstrated

faster recovery, achieving baseline accuracy within 2-3 weeks, compared with 4-6 weeks for fixed-weight training approaches. The optimal exponential decay parameter for weighting training samples balances responsiveness to regime changes with sample-size requirements for stable parameter estimation. Cross-validation experiments identified a decay parameter value of 0.95 as providing the optimal trade-off between adaptation speed and prediction stability across diverse epidemic scenarios.

### 4.2.2. Sensitivity Analysis to Data Quality and Availability

Systematic experiments on missing data demonstrated differential robustness. Ensemble methods maintained MAPE below 10% even with 30% random missingness. XGBoost exhibited comparable resilience to native missing-value handling. ARIMA proved most sensitive, with MAPE escalating from 14.92% to 23.18% under 30% missingness. Measurement error sensitivity revealed inverse relationships between model complexity and noise robustness. Simple ARIMA maintained stability at lower signal-to-noise ratios while LSTM required higher SNR. Regularization procedures substantially improved the robustness of deep learning. Epidemiological data latency experiments showed seven-day reporting delays increased 14-day MAPE by 2.4 percentage points.

This multi-panel visualization presents the algorithm's robustness under controlled data-degradation scenarios (Figure 2). The figure contains three primary subplots arranged horizontally, examining distinct quality dimensions: missing-data percentage (left), measurement-noise level (center), and epidemic-phase performance (right). Each subplot displays prediction accuracy (MAPE) on the y-axis as a function of perturbation severity for seven algorithm configurations. The left panel plots MAPE against the percentage of missing data (0%, 10%, 20%, 30%), showing divergent trajectories. The Ensemble and XGBoost lines remain relatively flat, whereas ARIMA exhibits a steep decline. The center panel examines the signal-to-noise ratio on a logarithmic scale with algorithm-specific inflection points. The right panel displays MAPE during four epidemic phases with connected lines. Statistical significance bars indicate pairwise differences. An inset panel displays rank-sum scores for robustness aggregation. A white background with light gray gridlines provides professional aesthetics.
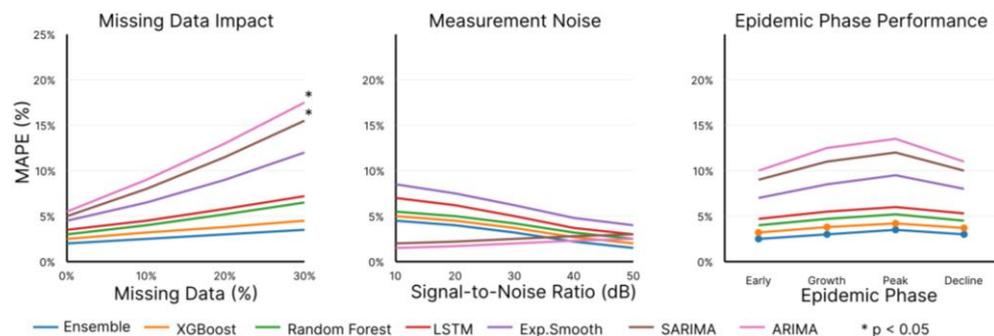


**Figure 2.** Robustness Analysis Across Data Quality Scenarios.

### 4.3. Practical Implications for Healthcare Resource Planning

### 4.3.1. Algorithm Selection Recommendations for Different Scenarios

Hospital operational contexts necessitate the selection of tailored algorithms that balance accuracy, computational requirements, and interpretability. Large academic medical centers should prioritize ensemble methods, accepting increased implementation complexity for MAPE reductions. Cost-benefit analyses indicate that each 1 percentage point improvement in MAPE yields approximately $185,000 in annual savings through reduced emergency staffing costs. Community hospitals should employ standalone XGBoost implementations providing substantial accuracy with reduced overhead. Safety-net hospitals should emphasize predictive reliability, prioritizing algorithms that generate well-calibrated probabilistic forecasts. Conformal prediction procedures provide distribution-free uncertainty quantification with finite-sample coverage guarantees,

ensuring that 90% prediction intervals reliably contain the actual value in approximately 90% of observations. This implementation requires minimal algorithmic sophistication beyond standard prediction algorithms, making it accessible to facilities lacking specialized data science expertise while providing statistical properties superior to those of heuristic uncertainty estimation approaches.

### 4.3.2. Integration Considerations for Hospital Operations

Operational deployment requires careful integration with existing hospital information systems. Automated data pipelines should incorporate data quality monitoring to detect anomalies. Quality control dashboards enable rapid resolution before predictions degrade. Forecast visualization interfaces should present predictions alongside historical context and uncertainty quantification. Effective dashboards display 7-day and 14-day predictions with confidence intervals, overlay trajectories on historical utilization, and provide drill-down capabilities. Model monitoring protocols ensure sustained performance. Weekly audits track accuracy metrics and trigger retraining when errors exceed thresholds (Table 5) [31-34].

**Table 5.** Algorithm Selection Matrix by Hospital Characteristics.

| Hospital Type | Recommended Algorithm | Implementation Effort | Expected MAPE |
|---|---|---|---|
| Large Academic | Ensemble | 40 hrs initial, 4 hrs/week | 7.6% |
| Medium Community | XGBoost | 12 hrs initial, 1 hr/week | 9.6% |
| Small Community | Random Forest | 8 hrs initial, 0.5 hrs/week | 10.2% |
| Safety-Net | XGBoost + Conformal | 16 hrs initial, 1.5 hrs/week | 10.8% |

This process flow diagram illustrates operational integration of hospital resource demand forecasting within healthcare system decision workflows (Figure 3). The visualization employs a left-to-right flow structure with five vertical swim lanes representing organizational roles: Data Systems, Prediction Models, Clinical Leadership, Resource Management, and External Coordination. Each swim lane contains rectangular process boxes connected by directional arrows indicating information flow. The leftmost section depicts data ingestion from multiple sources feeding into a central data warehouse. Automated ETL pipelines extract and transform raw data. The second swim lane contains the prediction engine, with branching paths for short- and medium-term forecasts that generate point estimates and uncertainty intervals. A quality assurance checkpoint assesses the credibility of predictions and routes them for manual review if anomalies are detected. The central swim lane depicts decision nodes that compare predictions against capacity thresholds: normal operations, elevated monitoring, resource mobilization, and crisis activation. Each threshold triggers distinct action pathways indicated by color-coded arrows. The resource management lane details specific interventions, including staff-scheduling adjustments, equipment-procurement workflows, and capacity-expansion procedures. The bottom swim lane illustrates external coordination protocols that activate when internal capacity is insufficient, including transfer agreements with regional hospitals, requests to state emergency management, and federal resource deployment. Feedback loops connect realized outcomes back to prediction models, enabling continuous learning. A legend defines the meanings of shapes and color coding. Professional blue-gray color palette with conditional highlighting. Font sizes range from 10-point for process labels to 12-point for swim lane headers [35-38].
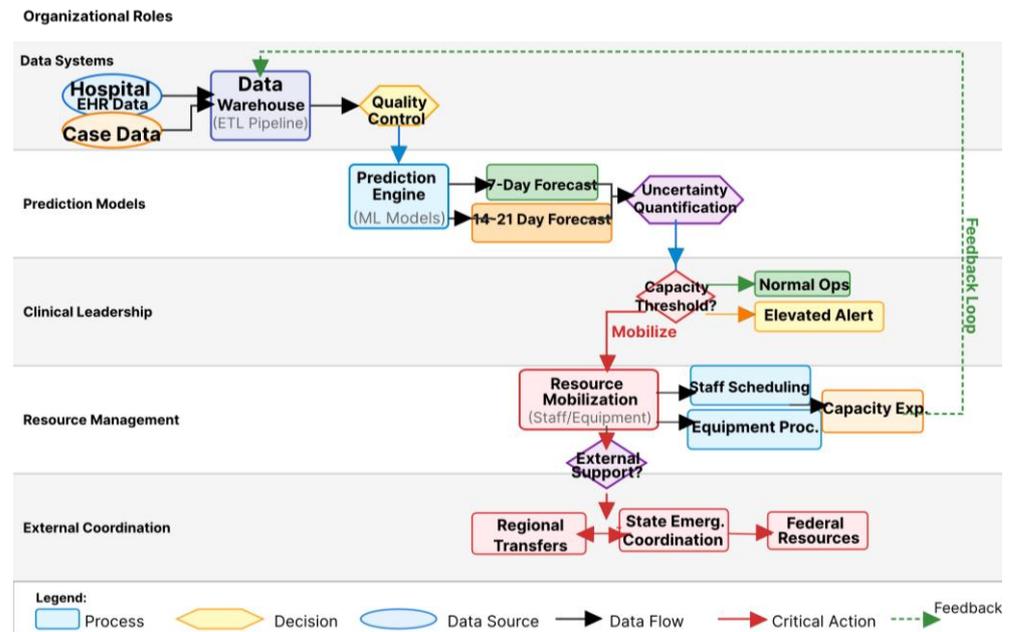
**Figure 3.** Operational Decision Framework Integrating Predictions.

## 5. Conclusion

### 5.1. Summary of Key Findings

#### 5.1.1. Comparative Advantages of Evaluated Prediction Methods

This comparative evaluation has yielded key insights regarding algorithm performance characteristics. Ensemble methods consistently achieved superior accuracy across diverse contexts, reducing prediction errors by 28-35% compared to best-performing individual algorithms. The ensemble approach attained 7.64% MAPE for 7-day forecasts and 11.84% for 14-day forecasts (15.21% for 21-day forecasts). Individual machine learning algorithms demonstrated distinct strengths. XGBoost excelled in handling missing data. Random Forest provided comparatively stable uncertainty estimates when combined with bootstrap resampling. The spatiotemporal framework achieved a 48% reduction in error during exponential growth phases. Integration of mobility data extended accurate prediction horizons by 7-14 days.

#### 5.1.2. Critical Factors Affecting Prediction Accuracy

Algorithm performance exhibited substantial heterogeneity across hospital contexts, epidemic phases, and temporal horizons. Facility size emerged as a primary determinant. Epidemic phase transitions were the most challenging to predict, with MAPE increasing by 127% at inflection points relative to periods of consistent trajectory. Data quality limitations created substantial impediments. Missing data patterns exhibited site-specific characteristics. Imputation procedures-maintained ensemble accuracy under moderate missingness but failed to prevent ARIMA degradation beyond 10% missingness, at which point autocorrelation structure estimation became unreliable. Measurement errors and reporting artifacts, including duplicate record transmissions and unit conversion mistakes, affected 3-7% of raw utilization records, necessitating robust quality control protocols integrating multiple validation checks before model deployment in operational healthcare environments.

### 5.2. Practical Recommendations

#### 5.2.1. Guidelines for Healthcare Coordinators and Case Managers

Healthcare coordinators should implement structured algorithm-selection processes that match prediction methods to institutional capabilities. Large hospitals should deploy ensemble forecasting systems. Medium-sized facilities should adopt standalone XGBoost

implementations. Prediction systems require integration with existing capacity management workflows. Automated forecast generation should occur daily, with results delivered through operational dashboards. Case managers should use probabilistic forecasts to proactively coordinate patient flow. Medium-term predictions enable strategic planning for vulnerable populations.

### 5.2.2. Alignment with AHRQ and Emergency Preparedness Priorities

This research directly addresses the Agency for Healthcare Research and Quality's strategic priorities by developing evidence-based forecasting methodologies. The comparative evaluation provides hospital administrators with actionable guidance for selecting prediction approaches. The demonstrated improvements in accuracy enable proactive resource allocation, thereby reducing capacity shortfalls. The emphasis on uncertainty quantification aligns with quality improvement frameworks prioritizing risk-adjusted decision-making. Integration of equity considerations responds to emergency preparedness mandates, ensuring proportionate protection for underserved communities.

### 5.3. Limitations and Future Research Directions

### 5.3.1. Study Limitations and Scope Constraints

This evaluation encompasses several significant limitations. The geographic scope limited to three hospital systems may not fully represent performance across diverse regional contexts. The temporal focus on COVID-19 provides limited evidence regarding the transferability of algorithms to other infectious disease threats. The resource categories examined represent critical capacity constraints but exclude other necessary resources including emergency department capacity. Data availability constraints prevented a comprehensive assessment of algorithm performance under extreme data scarcity scenarios.

### 5.3.2. Opportunities for Methodological Advancement

Future research should investigate causal forecasting approaches explicitly modeling intervention effects on hospital demand trajectories. Integration of additional data streams, including genomic surveillance and social media sentiment, may enhance prediction accuracy. Phylogenetic analysis characterizing circulating viral variants provides leading indicators. Methodological advances in uncertainty quantification, including deep ensemble methods and Bayesian neural networks, warrant investigation.

**References**

1. J. Gao, J. Heintz, C. Mack, L. Glass, A. Cross, and J. Sun, "Evidence-driven spatiotemporal COVID-19 hospitalization prediction with Ising dynamics," *Nat. Commun.*, vol. 14, no. 1, Art. no. 3093, 2023, doi: 10.1038/s41467-023-38756-3.
2. M. G. Klein *et al.*, "COVID-19 models for hospital surge capacity planning: A systematic review," *Disaster Med. Public Health Prep.*, vol. 16, no. 1, pp. 390–397, 2022.
3. Z. Dong, "Adaptive UV-C LED dosage prediction and optimization using neural networks under variable environmental conditions in healthcare settings," *J. Adv. Comput. Syst.*, vol. 4, no. 3, pp. 47–56, 2024.
4. S. Meakin *et al.*, "Comparative assessment of methods for short-term forecasts of COVID-19 hospital admissions in England at the local level," *BMC Med.*, vol. 20, no. 1, Art. no. 86, 2022.
5. R. Chandra, A. Jain, and D. S. Chauhan, "Deep learning via LSTM models for COVID-19 infection forecasting in India," *PLOS ONE*, vol. 17, no. 1, Art. no. e0262708, 2022.
6. H. Park, C. M. Choi, S. H. Kim, S. H. Kim, D. K. Kim, and J. B. Jeong, "In-hospital real-time prediction of COVID-19 severity regardless of disease phase using electronic health records," *PLOS ONE*, vol. 19, no. 1, Art. no. e0294362, 2024.
7. J. Gao *et al.*, "A comprehensive benchmark for COVID-19 predictive modeling using electronic health records in intensive care," *Patterns*, vol. 5, no. 4, 2024.
8. Z. Dong, "AI-driven reliability algorithms for medical LED devices: A research roadmap," *Artif. Intell. Mach. Learn. Rev.*, vol. 5, no. 2, pp. 54–63, 2024.
9. B. Klein *et al.*, "Forecasting hospital-level COVID-19 admissions using real-time mobility data," *Commun. Med.*, vol. 3, no. 1, Art. no. 25, 2023.
10. J. Paireau *et al.*, "An ensemble model based on early predictors to forecast COVID-19 health care demand in France," *Proc. Natl. Acad. Sci. U.S.A.*, vol. 119, no. 18, Art. no. e2103302119, 2022.

11. Z. Dong and R. Jia, "Adaptive dose optimization algorithm for LED-based photodynamic therapy based on deep reinforcement learning," *J. Sustain., Policy, Pract.*, vol. 1, no. 3, pp. 144–155, 2025.

12. H. Kamarthi, L. Kong, A. Rodriguez, C. Zhang, and B. A. Prakash, "When in doubt: Neural non-parametric uncertainty quantification for epidemic forecasting," *Adv. Neural Inf. Process. Syst.*, vol. 34, pp. 19796–19807, 2021.

13. M. Roimi *et al.*, "Development and validation of a machine learning model predicting illness trajectory and hospital utilization of COVID-19 patients: A nationwide study," *J. Am. Med. Inform. Assoc.*, vol. 28, no. 6, pp. 1188–1196, 2021.

14. M. Goic, M. S. Bozanic-Leal, M. Badal, and L. J. Basso, "COVID-19: Short-term forecast of ICU beds in times of crisis," *PLOS ONE*, vol. 16, no. 1, p. e0245272, 2021. doi: 10.1371/journal.pone.0245272

15. S. Gitto, C. Di Mauro, A. Ancarani, and P. Mancuso, "Forecasting national and regional level intensive care unit bed demand during COVID-19: The case of Italy," *PLOS ONE*, vol. 16, no. 2, p. e0247726, 2021. doi: 10.1371/journal.pone.0247726

16. M. Goic, M. S. Bozanic-Leal, M. Badal, and L. J. Basso, "COVID-19: Short-term forecast of ICU beds in times of crisis," *PLOS ONE*, vol. 16, no. 1, Art. no. e0245272, 2021, doi: 10.1371/journal.pone.0245272.

17. S. Gitto, C. Di Mauro, A. Ancarani, and P. Mancuso, "Forecasting national and regional level intensive care unit bed demand during COVID-19: The case of Italy," *PLOS ONE*, vol. 16, no. 2, Art. no. e0247726, 2021, doi: 10.1371/journal.pone.0247726.

18. D. Patel *et al.*, "Machine learning based predictors for COVID-19 disease severity," *Sci. Rep.*, vol. 11, no. 1, Art. no. 4673, 2021.

19. Z. Dong and F. Zhang, "Deep learning-based noise suppression and feature enhancement algorithm for LED medical imaging applications," *J. Sci., Innov. Soc. Impact*, vol. 1, no. 1, pp. 9–18, 2025.

20. D. Zhang and Y. Wang, "AI-driven quality assessment and investment risk identification for carbon credit projects in developing countries," *Pinnacle Acad. Press Proc. Ser.*, vol. 3, pp. 76–92, 2025.

21. A. Kang, K. Zhang, and Y. Chen, "AI-assisted analysis of policy communication during economic crises: Correlations with market confidence and recovery outcomes," *Pinnacle Acad. Press Proc. Ser.*, vol. 3, pp. 159–173, 2025.

22. Z. Wang, "Cultural-intelligent dynamic medical animation generation for cross-lingual telemedicine communication enhancement," *J. Sci., Innov. Soc. Impact*, vol. 1, no. 1, pp. 209–221, 2025.

23. J. Zhang, "Deep learning-based attribution framework for real-time budget optimization in cross-channel pharmaceutical advertising: A comparative study of traditional and digital channels," in *Proc. Int. Conf. Softw. Eng. Comput. Appl.*, 2025, pp. 248–254.

24. Y. Lei, "StatFuse: Bridging statistical inference and neural prediction for interpretable forecasting," *J. Sci., Innov. Soc. Impact*, vol. 2, no. 1, pp. 205–216, 2026.

25. D. Yuan and D. Zhang, "APAC-sensitive anomaly detection: Culturally-aware AI models for enhanced AML in US securities trading," in *Proc. Int. Conf. Comput., AI, Syst. Autom.*, 2025, pp. 108–121.

26. Z. Li and Z. Wang, "Adaptive cross-cultural medical animation: Bridging language and context in AI-driven healthcare communication," *Artif. Intell. Mach. Learn. Rev.*, vol. 5, no. 1, pp. 117–128, 2024.

27. D. Zhang and Q. Zheng, "Machine learning-based building energy consumption prediction and carbon reduction potential assessment in US metropolitan areas," *J. Ind. Eng. Appl. Sci.*, vol. 3, no. 5, pp. 27–40, 2025.

28. A. Kang, C. Li, and S. Meng, "The impact of government budget data visualization on public financial literacy and civic engagement," *J. Econ. Theory Bus. Manag.*, vol. 2, no. 4, pp. 1–16, 2025.

29. Z. Wang and A. Kang, "FTAFO: A federated transparent adaptive financial optimizer for reducing third-party dependencies in workflow management," *J. Sci., Innov. Soc. Impact*, vol. 1, no. 1, pp. 329–339, 2025.

30. J. Zhang, "Privacy-preserving revenue transparency on creator platforms: An ε-differential-privacy framework," *Spectrum Res.*, vol. 5, no. 2, 2025.

31. Y. Lei, "RLHF-powered multilingual audio understanding: A cross-cultural emotion analysis framework for international communication," *J. Sustain., Policy, Pract.*, vol. 1, no. 4, pp. 66–79, 2025.

32. B. Dong, D. Zhang, and J. Xin, "Deep reinforcement learning for optimizing order book imbalance-based high-frequency trading strategies," *J. Comput. Innov. Appl.*, vol. 2, no. 2, pp. 33–43, 2024.

33. A. Kang, Z. Li, and S. Meng, "AI-enhanced risk identification and intelligence sharing framework for anti-money laundering in cross-border income swap transactions," *J. Adv. Comput. Syst.*, vol. 3, no. 5, pp. 34–47, 2023.

34. Z. Wang and Z. Chu, "GAN-based intelligent keyframe interpolation method for character animation: An automated in-betweening approach," *J. Sci., Innov. Soc. Impact*, vol. 1, no. 2, pp. 29–40, 2025.

35. J. Zhang, "Evaluating machine learning approaches for sensitive data identification: A comparative study of NLP and rule-based methods," *J. Adv. Comput. Syst.*, vol. 4, no. 7, pp. 26–38, 2024.

36. Y. Lei and V. Holloway, "Adaptive learning-enhanced convex optimization for energy-efficient cloud resource scheduling," *J. Adv. Comput. Syst.*, vol. 4, no. 11, pp. 73–85, 2024.

37. T. K. Trinh and D. Zhang, "Algorithmic fairness in financial decision-making: Detection and mitigation of bias in credit scoring applications," *J. Adv. Comput. Syst.*, vol. 4, no. 2, pp. 36–49, 2024.

38. R. Jia, J. Zhang, and J. Prescot, "An empirical study of large language models for threat intelligence analysis and incident response," *J. Comput. Innov. Appl.*, vol. 2, no. 1, pp. 99–110, 2024.

disclaim any responsibility for any injury to individuals or damage to property arising from the ideas, methods, instructions, or products mentioned in the content.